

## WHAT IS CLAIMED IS

1                   1.     A method for reducing magnitudes of output traffic bursts in a  
2 streaming media cache comprises:  
3                   receiving a request from a first client system for a stream of media data, the  
4 stream of media data including a first streaming media data packet and a second streaming  
5 media data packet;  
6                   receiving a request from a second client system for the stream of media data;  
7                   receiving the first streaming media data packet from an upstream server, the  
8 first streaming media data packet including a delivery time;  
9                   determining a first modified delivery time for the first streaming media data  
10 packet;  
11                   determining a second modified delivery time for the first streaming media data  
12 packet, the first modified delivery time different from the second modified delivery time;  
13                   modifying the first streaming media data packet with the first modified  
14 delivery time to form a first modified first streaming media data packet;  
15                   modifying the first streaming media data packet with the second modified  
16 delivery time to form a second modified first streaming media data packet;  
17                   outputting the first modified first streaming media data packet to the first  
18 client system at the first modified delivery time; and  
19                   outputting the second modified first streaming media data packet to the second  
20 client system at the second modified delivery time.

1                   2.     The method of claim 1 wherein determining a first modified delivery  
2 time comprises adding a first delay value to the delivery time.

1                   3.     The method of claim 2 wherein the first delay value is selected from  
2 the range: 0 to approximately 500 milliseconds.

1                   4.     The method of claim 3 wherein the first delay value is pseudo-  
2 randomly selected from the range.

1                   5.     The method of claim 1 further comprising storing a payload portion of  
2 the first streaming media in a storage within the streaming media cache.

6. The method of claim 2 wherein the second streaming media data packet includes a delivery time, the method further comprising:

- determining a first modified delivery time for the second streaming media data packet;
- determining a second modified delivery time for the second streaming media data packet, the first modified delivery time different from the second modified delivery time;
- modifying the second streaming media data packet with the first modified delivery time to form a first modified second streaming media data packet;
- modifying the second streaming media data packet with the second modified delivery time to form a second modified second streaming media data packet;
- outputting the first modified second streaming media data packet to the first client system at the first modified delivery time; and
- outputting the second modified second streaming media data packet to the second client system at the second modified delivery time.

7. The method of claim 6 wherein determining the first modified delivery time for the second streaming media data packet comprises adding the first delay value to the delivery time of the second streaming media data packet.

8. The method of claim 1 further comprising:

- receiving a data file from the upstream server, the data file including a payload portion of the first streaming media data packet and a payload portion of the second streaming media data packet; and
- storing the data file in a storage within the streaming media cache.

9. A computer system for providing streaming media data to client systems with reduced magnitude traffic bursts comprises:

- a first thread configured to receive a request from a first client system and a second client system for a stream of data packets, wherein the stream includes a first data packet and a second data packet, the first thread also configured to specify a first client delay and a second client delay;
- a second thread coupled to the first thread, the second thread configured to receive the first data packet and the second data packet from an upstream server, wherein the first data packet specifies a first delivery time and the second data packet specifies a second

10 delivery time, the second thread also configured to form a first delayed first data packet in  
11 response to the first client delay and to form a second delayed first data packet in response to  
12 the second client delay, wherein the first delayed first data packet specifies a first delayed  
13 delivery time and the second delayed first data packet specifies a second delayed delivery  
14 time;

15 a third thread configured to receive the first delayed first data packet and to  
16 provide the first delayed first data packet to the first client system in response to the first  
17 delayed delivery time; and

18 a fourth thread configured to receive the second delayed first data packet and  
19 to provide the second delayed first data packet to the second client system in response to the  
20 second delayed delivery time.

10. The computer system of claim 9 wherein the second thread is  
configured to form the first delayed first data packet in response to the first client delay by  
adding the first client delay to the first delivery time.

11. The computer system of claim 10 wherein the first client delay is  
selected from the range: 0 to approximately 500 milliseconds.

12. The computer system of claim 9 further comprising a thread  
configured to store payload portions of the first data packet and payload portions of the  
second data packet in a memory

13. The computer system of claim 9 further  
wherein the second thread is also configured to form a first delayed second  
data packet in response to the first client delay and to form a second delayed second data  
packet in response to the second client delay, wherein the first delayed second data packet  
specifies a first delayed delivery time and the second delayed second data packet specifies a  
second delayed delivery time;

wherein the third thread is also configured to receive the first delayed second  
data packet and to provide the first delayed second data packet to the first client system in  
response to the first delayed delivery time specified therein; and

wherein the fourth thread is also configured to receive the second delayed  
second data packet and to provide the second delayed second data packet to the second client  
system in response to the second delayed delivery time specified therein.

1                   14.     The computer system of claim 13 wherein the second thread is also  
2 configured to form the first delayed second data packet in response to the first client delay by  
3 adding the second client delay to the second delivery time.

1                   15.     The computer system of claim 14 wherein the second client delay is  
2 pseudo-randomly selected.

1                   16.     A method for reducing peak output traffic bursts in a computer system  
2 where a first packet of data is scheduled to be delivered to more than one downstream client  
3 system substantially simultaneously comprises:

4                   delaying a packet delivery time for the first packet of data to be delivered to a  
5 first downstream client system; and

6                   delaying the packet delivery time for the first packet of data to be delivered to  
7 a second downstream client system;

8                   wherein the first packet of data is to be delivered to the first downstream client  
9 system at a time different than when the first packet of data is to be delivered to the second  
10 downstream client system.

11                   17.     The method of claim of claim 16 wherein the first packet of data is  
12 framed.

1                   18.     The method of claim 16 wherein the first packet of data comprises  
2 streaming media data.

1                   19.     The method of claim 16 wherein delaying the packet delivery time for  
2 the first packet of data to be delivered to the first downstream client system comprises  
3 delaying the first packet of data by a delay factor selected from 0-500 milliseconds.

1                   20.     The method of claim 19 further comprising delaying a packet delivery  
2 time for a second packet of data to be delivered to the first downstream client system by the  
3 delay factor.